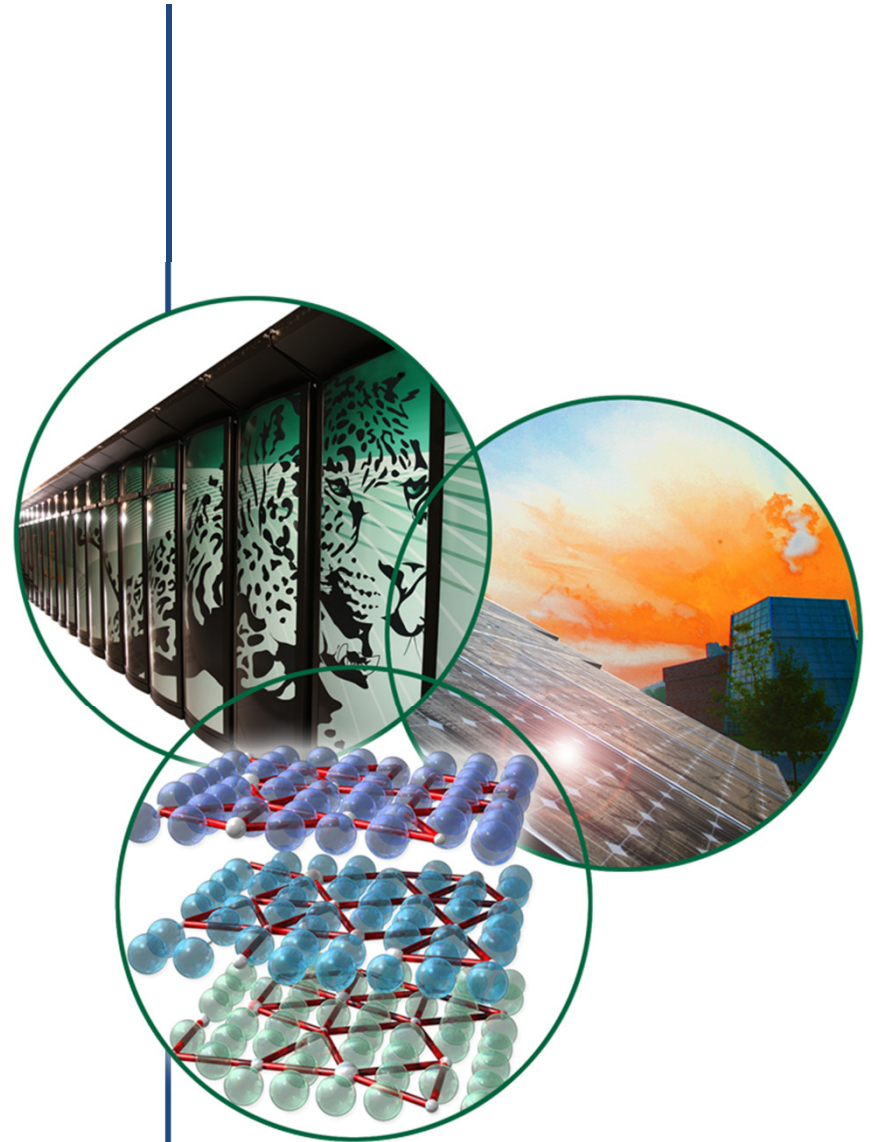# Toward Abstracting the Communication Intent in Applications to Improve Portability and Productivity

**Tiffany M. Mintz**, Oscar Hernandez, Christos Kartsaklis, David E. Bernholdt, Markus Eisenbach, Swaroop Pophale
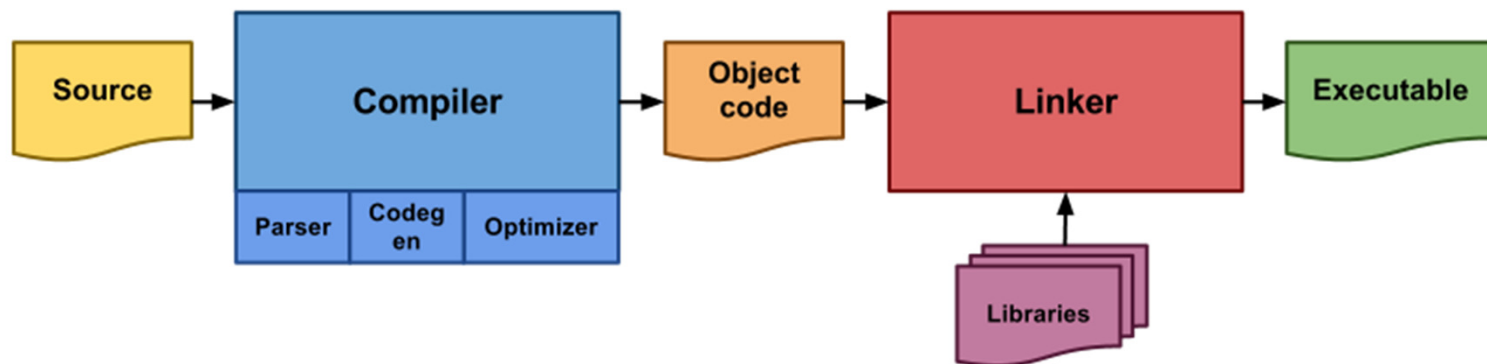
**U.S. DEPARTMENT OF ENERGY**

**OAK RIDGE NATIONAL LABORATORY**
MANAGED BY UT-BATTELLE FOR THE DEPARTMENT OF ENERGY

# Content

- **Motivation**

- **Related Work**

- **Approach**

- **Design & Implementation**

- **Experimental Results**

- **Next Steps**

Managed by UT-Battelle
for the U.S. Department of Energy

# Why This Matters

- **Communication missing from compiler's static analysis**
  - **Opportunities for automatic static optimization lost**

- **High level abstractions:**
  - **provide greater portability**
  - **enhances productivity**
  - **Easier to maintain**

OAK RIDGE
National Laboratory

# What Others Have Done: Bamboo

```
   6 #pragma bamboo olap
   7   for (it=0; it<num_iterations; it++){
   8    #pragma bamboo send{
   9       pack boundary values to message Buffer
  10       MPI_Isend(SendGhostcells) to left/right/up/down
  11    }
  12    #pragma bamboo receive{
  13       MPI_Recv(RecvGhostcells) from left/right/up/down
  14       unpack incoming data to ghost cells
  15    }
  16    MPI_Waitall();
  17    for(j=1; j < N/numprocs_Y -2; j++)
  18      for(i=1; i < N/numprocs_X -2; i++)
  19        V(j,i)= c*(U(j,i+1)+U(j,i-1)+U(j+1,i)+U(j-1,i))
  20    swap(U, V);
  21  }
  22 free U, V, SendGhostcells, RecvGhostcells
```

Olap-region

Send Blk

Recv Blk

Compute Blk

- **Annotations for MPI library calls**

- **Provides mechanism for static analysis and communication/computation overlap**

OAK RIDGE
National Laboratory

# What Others Have Done: OpenMPI (not the library implementation)

- **OpenMP-like directives for incremental parallelization**

- **Express collective communication**

```
#pragma ompi sync_sleeve
    for(i = 1; i <= YSIZE; i++)
#pragma ompi for
    for(j = 1; j <= XSIZE; j++)
        nu[i][j] = (u[i-1][j] + u[i+1][j]
                    + u[i][j-1] + u[i][j+1]) / 4.0;


#pragma ompi for reduction(+:res2)
    for(j = 1; j <= XSIZE; j++){
        tmp = (nu[i][j] - u[i][j]) / u[i][j];
        res2 += tmp * tmp;
    }
#pragma ompi single
    fprintf(stderr, "itr=%d res=%g\n", itr, res1);
```

OAK RIDGE
National Laboratory

# Our Solution

- **Use directive language extensions to assert communication**

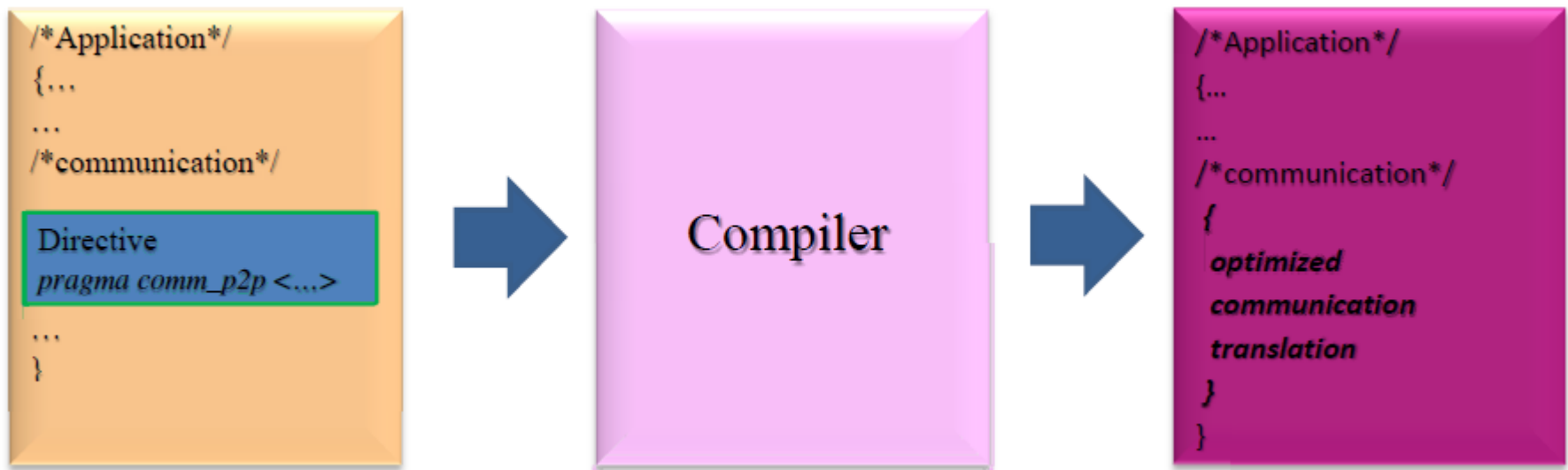- **Translate directives at compile time to communication calls**



**Assertion:**

A complete **sentence** which expresses the point.
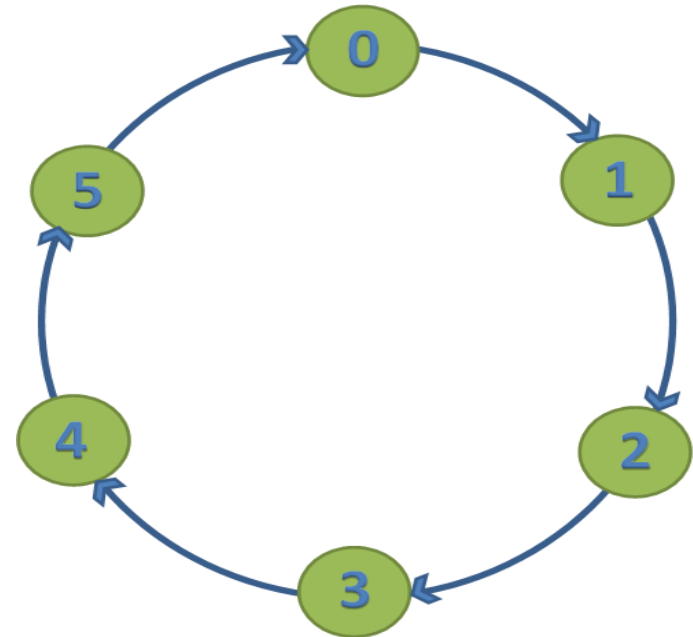
# Asserting Communication

- **Directives: comm_p2p, comm_parameters**

- **Clauses:**
  - **sender, receiver, sbuf, rbuf**
  - **sendwhen, receivewhen, count, place_sync, max_comm_iter, target**



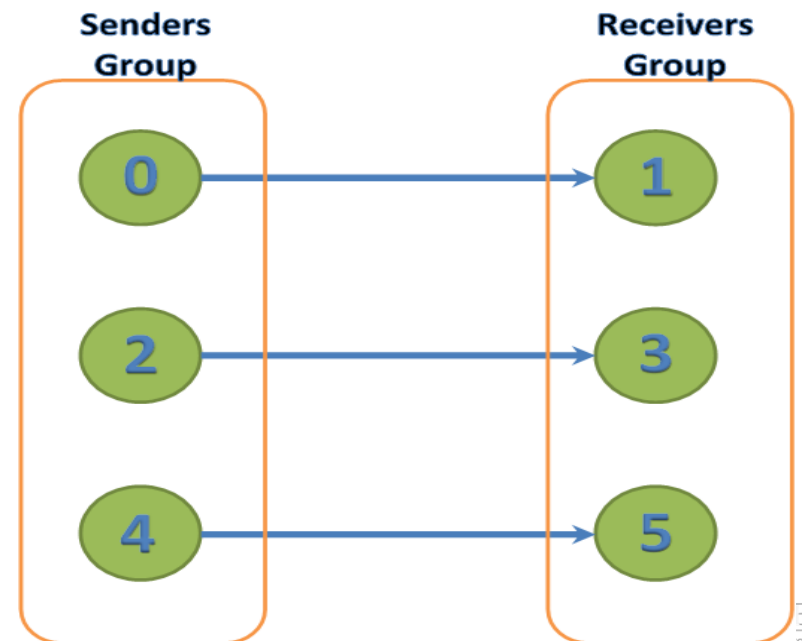Managed by UT-Battelle
for the U.S. Department of Energy

# Examples

Ring communication pattern

```
prev = (rank-1+nprocs)%nprocs;
next = (rank+1)%nprocs;
#pragma comm_p2p sender(prev) receiver(next) \
  sbuf(buf1) rbuf(buf2)
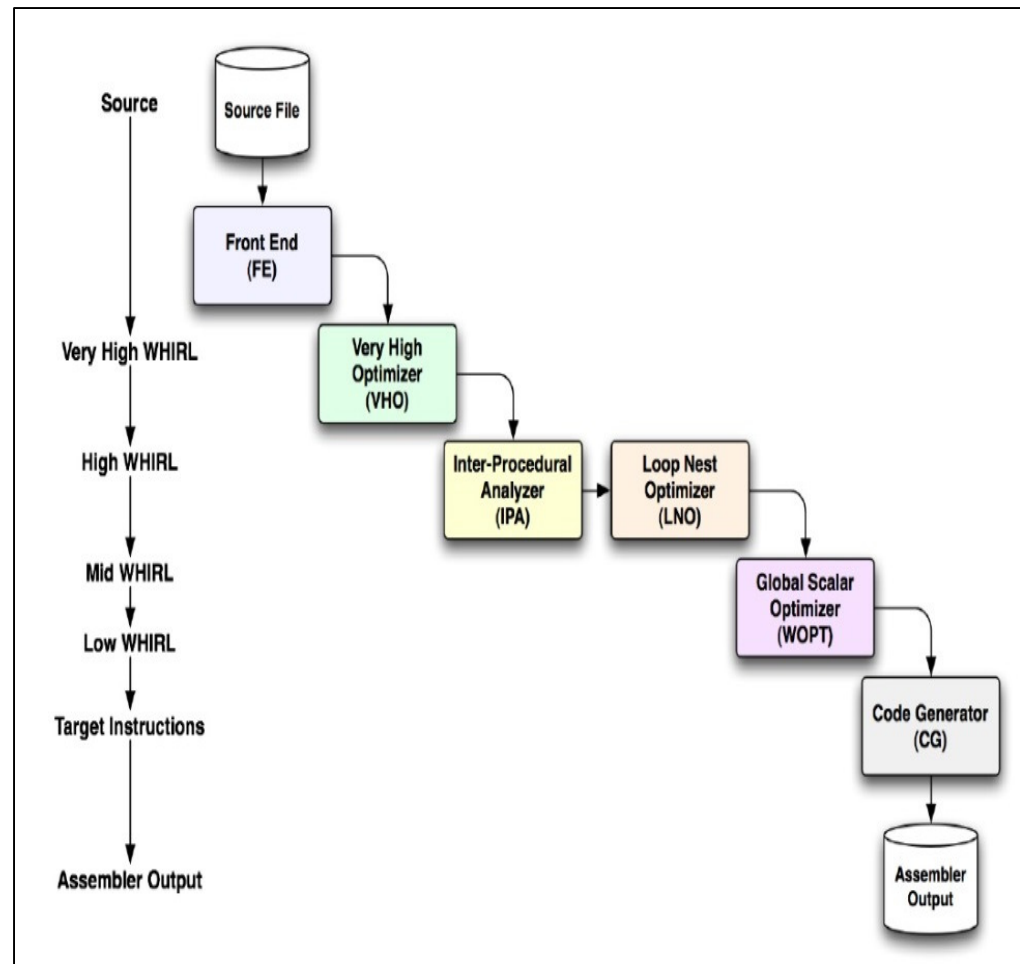```

Communication scoping and parameter inheritance

```
#pragma comm_parameters sender(rank-1) \
  receiver(rank+1) sendwhen(rank%2==0) \
  receivewhen(rank%2==1) count(size) \
  max_comm_iter(n) place_sync(END_PARAM_REGION)
{

for(p=0; p < n; p++)
  #pragma comm_p2p sbuf(&buf1[p]) rbuf(&buf2[p])
}
```

# Static Analysis and Optimizations

- **Communication scoping**

- **Communication/computation overlap**

- **Flexible Implementation**

- **Data type handling**
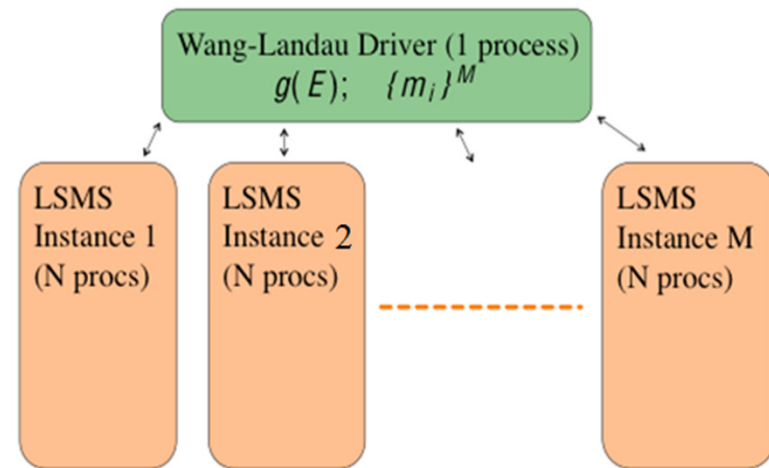
- **Synchronization reduction**



Open64 Compilation Phases

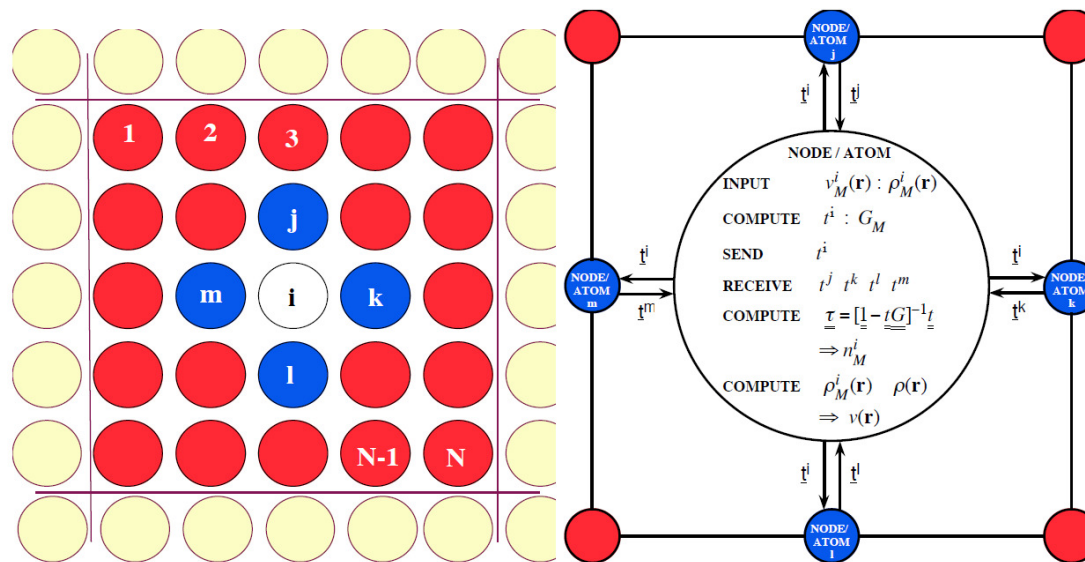# Scientific Application

## WL-LSMS

- ## Wang-Landau (WL)

  - Monte-Carlo calculation

- ## Locally Self-Consistent Multiple Scattering (LSMS)

  - First principles electronic structure calculation



**Wang-Landau Driver (1 process)**
$$g(E); \quad \{m_i\}^M$$

LSMS Instance 1 (N procs)

LSMS Instance 2 (N procs)

LSMS Instance M (N procs)

**WL-LSMS Organizational View**

OAK RIDGE National Laboratory

# WL-LSMS Communication

- ## Local Interaction Zone (LIZ)
  - ### Within each LSMS instance
  - ### Master-Worker process topology
  - ### Point-to-point communication



**LIZ Communication Pattern**

# Single Atom Data Communication

```
if(comm.rank==from)
{
    int pos=0;
    MPI_Pack(&local_id,1,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.jmt,1,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.jws,1,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.xstart,1,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.rmt,1,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(atom.header,80,MPI_CHAR,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.alat,1,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.efermi,1,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.vdif,1,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.ztotss,1,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.zcorss,1,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(atom.evec,3,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.nspin,1,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.numc,1,MPI_INT,buf,s,&pos,comm.comm);

    t=atom.vr.n_row();

    MPI_Pack(&t,1,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.vr(0,0),2*t,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.rhotot(0,0),2*t,MPI_DOUBLE,buf,s,&pos,comm.comm);

    t=atom.ec.n_row();

    MPI_Pack(&t,1,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.ec(0,0),2*t,MPI_DOUBLE,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.nc(0,0),2*t,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.lc(0,0),2*t,MPI_INT,buf,s,&pos,comm.comm);
    MPI_Pack(&atom.kc(0,0),2*t,MPI_INT,buf,s,&pos,comm.comm);

    MPI_Send(buf,s,MPI_PACKED,to,0,comm.comm);
}
if(comm.rank==to)
{
    MPI_Status status;
    MPI_Recv(buf,s,MPI_PACKED,from,0,comm.comm,&status);

    int pos=0;
    MPI_Unpack(buf,s,&pos,&local_id,1,MPI_INT,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.jmt,1,MPI_INT,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.jws,1,MPI_INT,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.xstart,1,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.rmt,1,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,atom.header,80,MPI_CHAR,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.alat,1,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.efermi,1,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.vdif,1,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.ztotss,1,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.zcorss,1,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,atom.evec,3,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.nspin,1,MPI_INT,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.numc,1,MPI_INT,comm.comm);

    MPI_Unpack(buf,s,&pos,&t,1,MPI_INT,comm.comm);

    if(t<atom.vr.n_row())
        atom.resizePotential(t+50);

    MPI_Unpack(buf,s,&pos,&atom.vr(0,0),2*t,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.rhotot(0,0),2*t,MPI_DOUBLE,comm.comm);

    MPI_Unpack(buf,s,&pos,&t,1,MPI_INT,comm.comm);

    if(t<atom.nc.n_row())
        atom.resizeCore(t);

    MPI_Unpack(buf,s,&pos,&atom.ec(0,0),2*t,MPI_DOUBLE,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.nc(0,0),2*t,MPI_INT,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.lc(0,0),2*t,MPI_INT,comm.comm);
    MPI_Unpack(buf,s,&pos,&atom.kc(0,0),2*t,MPI_INT,comm.comm);
}
```

← Original communication source code
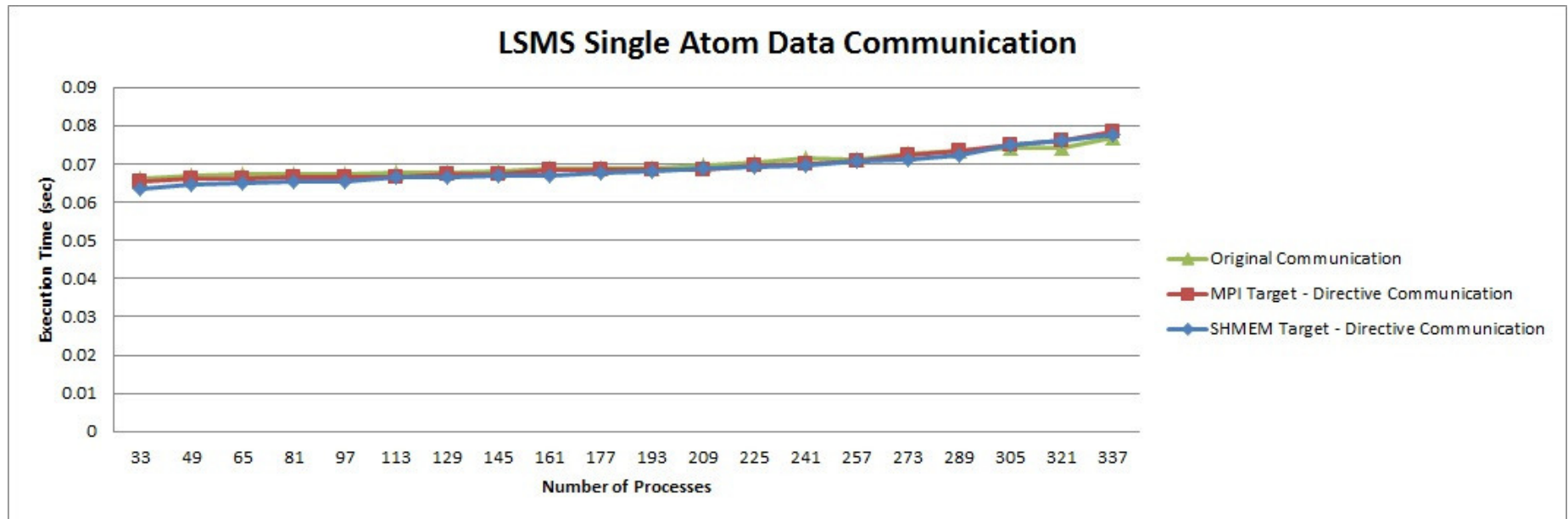
Communication using directives

```
#pragma comm_parameters sendwhen(rank==from_rank)\
  receivewhen(rank==to_rank) \
  sender(from_rank) receiver(to_rank)
{
    #pragma comm_p2p sbuf(scalaratomdata) \
      rbuf(scalaratomdata) count(1)
    { }

    #pragma comm_p2p  sbuf(vr,rhotot) \
      rbuf(vr,rhotot) count(size1)
    { }

    #pragma comm_p2p sbuf(ec,nc,lc,kc) \
      rbuf(ec,nc,lc,kc) count(size2)
    { }
}
```

# Performance Comparison



LSMS Single Atom Data Communication

- **Experiments using MPI and SHMEM translations**

- **Performance comparable to original source code**

# Spin Configurations Communication

Original communication source code ⟶

```
if(rank==0)
{
  for(int p=0; p<num_types; p++)
  {
    if(n==0)
      /*write to local space*/
    else
      MPI_Isend(&ev[3*p],3,MPI_DOUBLE,n,l,comm.
        comm,&request[n_req++]);
  }
  for(int i=0; i<n_req; i++)
    MPI_Wait(&request[i],&status);
} else {
  for(int p=0; p<num_local; p++)
  {
    MPI_Irecv(&local.atom[p].evec[0],3,MPI_DOUBLE
      ,0,p,comm.comm,&request[p]);
  }
  for(int i=0; i<num_local; i++)
    MPI_Wait(&request[i],&status);
}
```
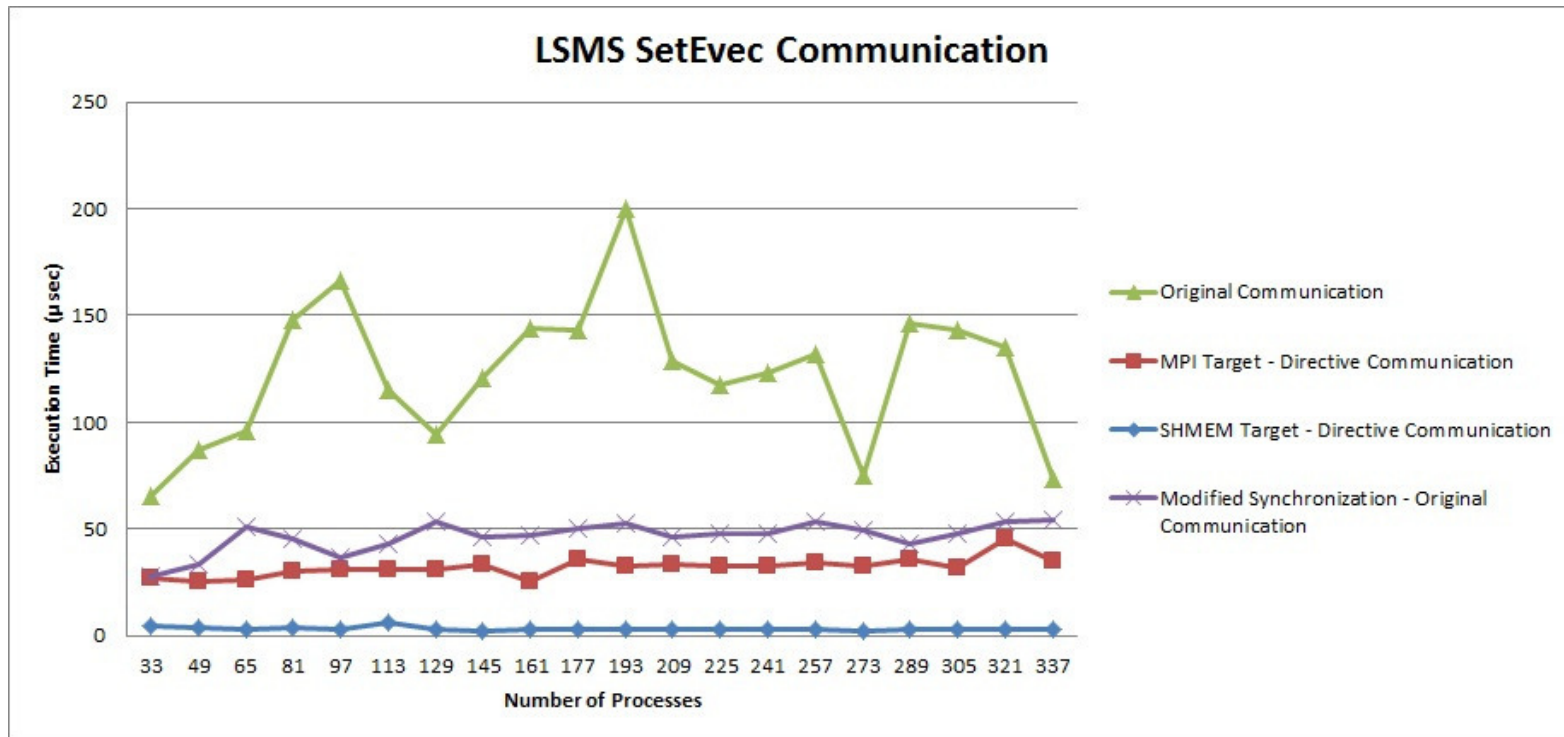
```
while((rank == 0 && send_p < num_types)
    || (rank != 0 && recv_p < num_local))
{
 if(n==0)
   /*write to local space*/
 else
 {
   #pragma comm_parameters sendwhen(rank == 0)\
     receivewhen(rank != 0)  sender(rank0) \
     receiver(rcv_rank)  count(3) \
     max_comm_iter(num_types) \
     place_sync(END_PARAM_REGION)
   {
    while((rank == 0 && n == types[send_p].node)
        || (rank != 0 && recv_p < num_local))
    {
     #pragma comm_p2p sbuf(&ev[3*send_p]) \
       rbuf(&local.atom[p].evec[0])
     {
       calculateCoreState(comm,lsms,local,recv_p
          ,!core_states_done);
     }
    }
   }
 }
}
```

⟵ Communication/computation overlap using directives

OAK RIDGE National Laboratory

# Communication Comparison



**LSMS SetEvec Communication**

- Original Communication
- MPI Target - Directive Communication
- SHMEM Target - Directive Communication
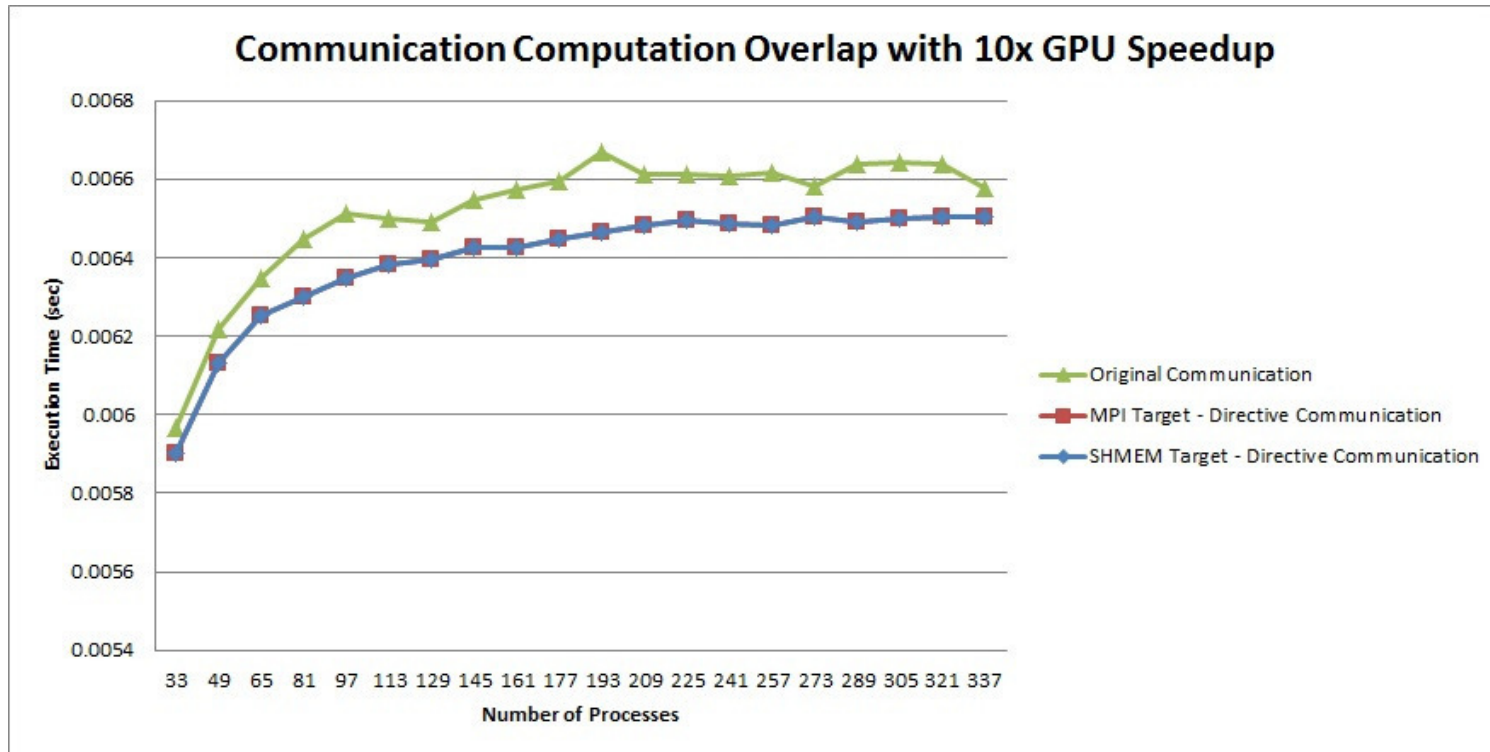- Modified Synchronization - Original Communication

- **Original synchronization caused poor performance**

- **After modifying synchronization:**
  - **MPI translation 1.4x speedup**
  - **SHMEM translation 14.5x speedup**

Managed by UT-Battelle
for the U.S. Department of Energy

# Communication/Computation Overlap



- **Current Computation/Communication ratio: 19 to 1**

- **Estimate 10x speedup with GPU acceleration**

OAK
RIDGE
National Laboratory

# Looking Forward

- **Summary**
  - Higher abstraction for message passing communication
  - Communication aware compiler
  - Static analysis and optimization for message passing

- **What's Next**
  - Develop assertions for many-to-one, one-to-many patterns
  - Extend data flow analysis
  - Implement cost model for automated selection of communication calls

OAK
RIDGE
National Laboratory

# QUESTIONS?

OAK
RIDGE
National Laboratory